

米国のがん統計に用いられている 数理モデルの概観

片野田 耕太[†]

(受付 2010 年 12 月 27 日; 改訂 2011 年 4 月 28 日; 採択 5 月 24 日)

要 旨

がんの統計情報は、国のがん対策の立案と評価のために不可欠である。わが国では、死亡統計が約 1 年遅れで最新データが公表されるのに対して、罹患統計は、最新データが公表されるのは約 5 年遅れであり、罹患統計の遅れの解消が大きな課題となっている。米国では、2010 年現在、死亡統計では状態空間モデル、罹患統計では空間モデル・時空間モデル・Joinpoint 回帰モデルの組合せにより、それぞれ 3 年後、4 年後の予測が行われ、いずれもリアルタイムの推計値が公表され政策利用されている。わが国においても、短期予測による最新の罹患統計を整備し、科学的根拠に基づくがん対策の情報基盤を構築してゆく必要がある。

キーワード：がん，死亡，罹患，予測，状態空間モデル，時空間モデル。

1. 緒言

がんの統計情報は、国のがん対策の立案と評価のために不可欠である。限られた資源をがん対策のために効果的かつ効率的に利用するためには、科学的証拠に基づいた政策決定が必要であり、がんの統計情報による実態把握および評価 (monitoring and evaluation) は、がん対策の立案と評価の根幹を成す部分である (World Health Organization, 2002)。がん対策に用いられるがんの統計情報には、死亡率、罹患率、生存率といった疾病統計から、予防・危険因子の動向、がん検診の動向、さらには患者の生活の質に至るまで、様々なものがある。これらの中で、死亡率および罹患率は、対象集団全体の疾病動向を示すという意味で重要な指標である。がん対策には、がん検診や診断・治療技術の普及のように、疾病統計に短期的に反映される政策課題も多い。疾病動向をリアルタイムで把握することで、より機動的な政策の立案と評価が可能となる。

死亡率および罹患率は、人口動態統計や地域がん登録などの疾病把握制度に基づいており、最新のデータが利用可能になるまでに一定の期間が必要となる。例えば、わが国における死亡統計は、人口動態統計に基づいており、全国規模の最新データが公表されるのは約 1 年遅れである (例えば、2009 年死亡データは 2010 年秋に公表される)。一方、がんの罹患統計は、地域がん登録に基づいており、独立行政法人国立がん研究センターがん対策情報センター (以下、がん対策情報センター) から全国推計値の最新データが公表されるのは約 5 年遅れである (例えば、2005 年罹患データが 2010 年に公表される) (Matsuda et al., 2011)。最新のデータがどの程

[†] 独立行政法人国立がん研究センター がん対策情報センターがん統計研究部：〒104-0045 東京都中央区築地 5-1-1；kkatanod@ncc.go.jp

度早く利用可能になるかは、がん統計の分野では「即時性(timeliness)」と呼ばれている。わが国の全国規模のがん統計情報は、死亡データではほぼ即時性が実現されているが、罹患データでは即時性の向上が課題である。

米国では、がん統計データの即時性を向上させるために、入手可能な最新データに数理モデルを適用して、数年後の予測が行われている。その結果、米国対がん協会(American Cancer Society, 以下 ACS)から毎年発表されるがん統計報告では、がん死亡およびがん罹患について、その年のリアルタイムの推計値が掲載される。例えば2010年の場合、CA: A Cancer Journal for Clinicians 誌に“Cancer Statistics, 2010”という標題の論文が掲載され、2007年までの死亡データ、および2006年までの罹患データに基づいて、いずれも2010年の推計値が発表されている(Jemal et al., 2010)。このような米国の例は、わが国のがん統計において、特に罹患データの即時性を向上させる上で参考になると思われる。本稿では、米国においてがん死亡およびがん罹患の短期予測に用いられている数理モデルを概観し、今後のわが国のがん統計への示唆を得ることを目的とする。

2. 短期予測の手法の変遷

2.1 概観

表1に、米国において用いられてきた短期予測の手法の概要を示す。死亡統計の2003年推計まで、および罹患統計の2006年推計までは、二次回帰モデルが共通で用いられてきた(Wingo et al., 1998)。その後、死亡統計については2004年推計から状態空間モデルが用いられ(Tiwari et al., 2004)、罹患統計については2007年推計から、空間モデル・時空間モデル・Joinpoint 回帰モデルの3つのモデルを組み合わせた手法が用いられるようになった(Pickle et al., 2007)。死亡、罹患とも、現在(2010年時点)でもこれらの手法が踏襲されている。

2.2 旧モデル(二次回帰モデル)

死亡の2003年推計まで、および罹患の2006年推計まで、短期予測には二次回帰モデルが用いられていた。この手法は、死亡数または罹患数に二次回帰モデルをあてはめ、残差に自己回帰過程を適用するものである(次式)(Wingo et al., 1998)。

$$R_t = \beta_0 + \beta_1 t + \beta_2 t^2 + u_t$$

$$u_t = \gamma_1 u_{t-1} + \cdots + \gamma_p u_{t-p} + e_t$$

(R_t : 時点 t の年齢調整罹患率, u_t : 残差, p : 最大ラグ,

$\beta_i (i=0, 1, 2)$ および $\gamma_i (i=1, \dots, p)$: 各項の係数, e_t : 誤差項)

死亡推計においては、最終的な予測値は、モデルによる予測値(点推定値)、その信頼区間の上下限、およびそれらの上下限と点推定値との中点、の5つから、主観的判断で選択されていた(Tiwari et al., 2004)。また、罹患推計においては、全米の罹患率は、全米人口の10%程度を占める一部地域の罹患率で代表されとの前提を置いていた。死亡統計における2004年の手法変更、および罹患統計における2007年の手法変更は、これらの点を改善することが主眼となった。

2.3 死亡推計の新技术法(状態空間モデル)

死亡統計については、2004年推計から、状態空間モデルが用いられている。この手法は、死亡数に線形回帰モデルをあてはめ(measurement equation)、その回帰係数に一階の自己回帰線形モデルをあてはめる(transition equation)ものである(次式)。

表 1. 米国のがん死亡・罹患統計に用いられてきた短期予測の手法.

死亡/罹患	使用された年 ^a	モデル	概要	使用データ	推計対象年	文献
死亡	1999年～2003年	二次回帰モデル	死亡数に二次回帰モデルをあてはめ、残差に自己回帰過程を適用する。	全米健康統計センター (National Center for Health Statistics) の 1969～X年 ^b データ	(X+3)年 ^b	Wingo et al., 1998
	2004年～現在 (2010年)	状態空間モデル	死亡数に時間依存性回帰係数をもつ線形回帰モデルをあてはめ、回帰係数に一階の線形自己回帰モデルを適用する。	全米健康統計センターの1969年～X年 ^b 死亡データ	(X+3)年 ^b	Tiwari et al., 2004
罹患	1999年～2006年	二次回帰モデル	罹患数 ^d に二次回帰モデルをあてはめ、残差に自己回帰過程を適用する (死亡の1999年～2003年の手法と同じ)。2004年～2006年推計にはdelay adjust ^e が適用されている。	SEER の1979年～X年 ^c データ (9地域) ^f	(X+4)年 ^c	Wingo et al., 1998
	2007年～現在 (2010年)	空間モデル・時空間モデル・Joinpoint回帰モデル	空間モデル、時空間モデル、Joinpoint回帰モデルの3段階から成る。空間モデルは、county毎・年齢階級毎の罹患数に、死亡率、人口学的・生活習慣共変量などを説明変数とするポワソン回帰をあてはめる。時空間モデルは、空間モデルで推計されたcounty別の罹患数に、時間自己相関を加味した二次回帰モデルをあてはめる。時空間モデルにより推計されたcounty毎の罹患数を合計して全国値を求め、全国値にdelay-adjust ^e を行う。Joinpoint回帰モデルは、上記で求められた全国値の時系列 (外数) に対して折れ線をあてはめ、直近の線分外挿することによって予測を行う。	罹患データ: 北米中央がん登録室協議会 (NAACCR) の1995年～X年 ^c 罹患データ (44州) 死亡データ: 全米健康統計センターの死亡データ 共変量データ: Behavioral Risk Factor Surveillance System (BRFSS) などの調査データ	(X+4)年 ^c	Pickle et al., 2007

a. 全米対がん協会 (American Cancer Society) が当該手法を用いて推計した年。
b. 死亡は3年後の予測を行う (例: 1999年～2007年死亡のデータを用いて2010年の推計をする)。ただし、年によっては使用データの年および推計対象年が異なる場合がある。
c. 罹患は4年後の予測を行う (例: 1995年～2006年罹患のデータを用いて2010年の推計をする)。
d. 罹患数は、SEERの対象とする9地域の年齢階級別罹患率に全国年齢階級別人口を乗じたものを合計して求める。
e. 集計機関に報告される罹患数には遅れと誤りがあることが知られており、それを定量的に補正する手法 (delay-adjustと呼ばれる) が開発され、2004年罹患推計以降利用されている (Olegg et al., 2002)。
f. SEERの9地域: Atlanta, Connecticut, Detroit, Hawaii, Iowa, New Mexico, San Francisco-Oakland, Seattle-Puget Sound, およびUtah

$$D_t = F_t \theta_t + e_t$$

$$\theta_t = G_t \theta_{t-1} + \eta_t$$

(D_t : 時点 t の死亡数, F_t および G_t : 係数行列,

θ_t : 未知パラメータベクトル, e_t : 誤差項, η_t : 誤差ベクトル)

状態空間モデルによる新手法の妥当性を検証した論文では、肺がん、前立腺がん、リンパ腫、女性の乳がんなど、多くのがん種で新手法の方が旧手法(二次回帰モデル+主観的判断)よりも予測精度がよかった、つまり実測値との乖離が小さかったと報告している(Tiwari et al., 2004)。ただ、新手法は、複数の短期的な変化がある場合に適合度がよい反面、短期的な変動の影響を受けやすく、予測値が安定しない傾向がある。例えば女性の大腸がんでは従来型の二次回帰モデルの方が予測値が安定していた(Tiwari et al., 2004)。

州レベルのがん死亡数の推計は、状態空間モデルによる推計を州別に行った後、州レベルの推計値の合計と全国レベルの推計値との差を各州に比例配分する。州レベルの推計では、状態空間モデルはランダムな変動を拾い過ぎる結果、旧モデル(二次回帰モデル)より予測精度が低かった(Tiwari et al., 2004)。政策利用を目的とした手法開発においては、疾病や対象集団ごとに予測精度が最も高い手法を別々に選択すべきであるという考え方と、国の政策決定に用いる手法は一つに統一すべきであるという考え方がある。米国では後者が重視され、死亡推計には空間状態モデルを統一的に用いるという決定がなされた。

2.4 罹患推計の新手法(3つのモデルの組合せ)

罹患データについては、2007年推計から3つのモデルを組み合わせた手法が用いられている(Pickle et al., 2007)。この手法は、空間モデル、時空間モデル、およびJoinpoint回帰モデルの3段階から成る。最初の空間モデルは、郡(county)ごとの各年罹患数をポワソン回帰で推計する(Pickle et al., 2003)。二番目の時空間モデルは、最初の空間モデルで推計された郡ごとの罹患数に、時間的空間的自己相関を考慮した二次回帰モデルをあてはめる。時空間モデルで推計された郡ごとの罹患数を合計することで、全国の罹患数を求める。最後のJoinpoint回帰モデルは、時系列に折れ線をあてはめる手法であり(Kim et al., 2000)、全国値の時系列に折れ線をあてはめ、直近の線分を外挿することで予測を行う。

最初の空間モデルは、罹患データがない(つまり地域がん登録が整備されていない)地域を含めて全米の郡別の罹患地図を作成するために開発されたものである(Pickle et al., 2003)。郡ごと・年齢階級ごとの罹患率に、年齢階級中央値、死亡率、人口学的共変量、生活習慣共変量を説明変数とするポワソン回帰をあてはめる(次式)。

$$\ln(\lambda_{i,j}) = f(a_j)\beta + \ln(m_{i,j})\gamma + X_i'\delta + Y_i'\zeta + e_s$$

($\lambda_{i,j}$: 郡 i , 年齢階級 j の罹患率, $f(a_j)$: 年齢階級 j の年齢中央値 a_j の三次関数,

$m_{i,j}$: 郡 i , 年齢階級 j の死亡率, X_i : 人口学的共変量ベクトル,

Y_i : 生活習慣共変量ベクトル, β, γ, δ および ζ : 各項の係数, e_s : 郡 i が属する地域 s の誤差項)

人口学的あるいは生活習慣に関する共変量には、郡ごとあるいは医療圏(health service area)ごとの収入、教育、住居、人種分布、都会・田舎の別、医療・がん検診機関の利用可能性、医療保険カバー率、喫煙率、肥満、がん検診受診率、および死亡率が含まれる。これらのデータは、米国疾病対策センター(CDC)、保健福祉省(Department of Health and Human Services)が運営するデータベースで入手されている(Health Resources and Services Administration, U.S. Department of Health and Human Services, 2010; National Center for Chronic Disease Prevention and Health Promotion, 2010; Center for Disease Control and Prevention, 2010)。

二番目の時空間モデルは、最初の空間モデルで推計された郡別の罹患数に、時空間自己相関を加味した二次回帰モデルをあてはめる。この二次回帰モデルでは、地域ごとに異なる時間効果と、時間的空間的自己相関をモデル化している。さらに、ポワソン分布に基づいて期待される変動からの超過を過分散項としてモデル化している。これにより推計された郡ごとの罹患数を合計して全国値を求める。なお、地域がん登録によって把握される罹患数には、遅れと誤りがあることが知られており、それを数理的に補正する手法(遅延補正(delay-adjust)と呼ばれる)が開発されている(Clegg et al., 2002)。2004年の推計以降、モデルによって推計された全国値に対して、この遅延補正が適用されている。

最後の Joinpoint 回帰モデルは、上記で求められた全国値の時系列(対数)に対して、変曲点を求め、各線分に対数線形モデルをあてはめる(Kim et al., 2000)。直近の線分を外挿することで4年後の予測値を得る。Joinpoint 回帰モデルは、本来予測を目的としたものではなく、死亡率や罹患率の増減およびその変化の判定を目的とした手法であるが、二次回帰モデルや空間状態モデルと比べて、全国レベル、州レベルともに予測精度がよい、という判断で罹患の短期予測に用いられている(Pickle et al., 2007)。

3つのモデルの組合せによる新手法を旧手法(二次回帰モデル)と方法論的に比較すると、旧手法が一部地域(National Cancer Institute (NCI)の Surveillance Epidemiology and End Results (SEER)にデータ提出をしている9地域)の平均罹患率を代表値として全国に適用していたのに対して、新手法は北米中央がん登録室協議会(NAACCR)の44州(2009年推計までは41州)の罹患率を利用している。その結果、罹患データの人口カバー率(データ提供地域の人口の、全国人口に占める割合)が旧手法で約10%であったのに対して、新手法では約90%である。これは、罹患データのソースである地域がん登録が、一定の精度を満たした形で全国的に整備されたことを反映している。SEERの9地域は社会経済指標が比較的高い都市部が多く、乳がんの罹患率が高い、全国に比べて喫煙率が低いという特徴があるため(Pickle et al., 2007)、この9地域で全国値を代表させることには問題があったが、新手法ではこれが改善された。また旧手法では、州レベルの推計を行う際に罹患/死亡比が全国一律であるという前提を置いていた。これは、がん患者の生存率に地域差がない、という前提と実質的に同じであり、非現実的な仮定である。新手法では地域別の罹患率を直接データとして使っているため、この問題が改善されている。さらに、空間推計に生活習慣や医療水準に関連する要因を加えている点、時空間モデルでは空間的、時間的自己相関を考慮している点なども、旧手法にはなかった新手法の特徴である。

新手法の妥当性を検証した論文では、空間推計と時間推計(予測)に分けて新手法と旧手法の精度を実測値との比較により評価している。全国レベルの空間推計に関しては、男性の肺がん、女性の乳がんなどでは新手法の方が、前立腺がん、女性の大腸がんなどでは旧手法の方が高い予測精度であった。州レベルでも性・がん種によって結果が異なり、単純な優劣をつけることは困難だが、新手法の方が実測値から大きく乖離することが少ない、という点が考慮されて、新手法が統一的な手法として採用された(Pickle et al., 2007)。時間推計(予測)に関しては、全国レベル、州レベルとも新手法(時空間モデルと Joinpoint 回帰モデルの組合せ)の精度が、多くの性・がん種の組合せで(全国レベルでは男性の肺・大腸・食道がん、女性の乳・皮膚がんなど、州レベルでは男女とも大腸・食道・皮膚がん、リンパ腫など)、旧手法より高かったため採用された(Pickle et al., 2007)。

3. わが国の現状と課題

3.1 現状

わが国では、前述の通り死亡統計は約1年遅れで実測の公表値が入手可能であり、人口動態

統計という全数調査に基づいているため地域別データの入手も可能である(1995年以降のデータは、毎年都道府県レベルまで公表されている)(総務省統計局, 2010)。一方、がんの罹患統計は、府県単位で行われている地域がん登録に基づいており、中央集計機関であるがん対策情報センターから最新のデータが公表されるのは約5年遅れとなっている。地域がん登録は2010年12月現在、38道府県市で実施されているが、がん対策情報センターが行う全国推計に含まれる地域は12府県、人口カバー率で約25%である(2005年推計)(独立行政法人国立がん研究センターがん対策情報センター, 2010a)。全国推計は、一定の精度基準を満たした地域がん登録を選んで行われる。全国推計に含まれる地域が10道府県しかないということは、38道府県のうち残りの三分の二はその基準を満たしていないことを意味している。都道府県別の罹患データについては、地域がん登録が実施されている道府県については「全国がん罹患モニタリング集計」としてがん対策情報センターから公表されているが(2010年12月現在、2005年罹患が最新)(Matsuda et al., 2011)、精度のばらつきが大きいため解釈には注意が必要である(独立行政法人国立がん研究センターがん対策情報センター, 2010b)。

がん死亡、罹患ともに、年次別・男女別・がん種別・年齢階級別の集計データががん対策情報センターで公開されている(独立行政法人国立がん研究センターがん対策情報センター, 2010c)。同サイトでは、生存率データ(7府県の地域がん登録に基づく)、都道府県別喫煙率、都道府県別がん検診受診率(いずれも国民生活基礎調査に基づく)なども公開されている。

3.2 課題

わが国のがん統計の最大の課題は、罹患の年次推移データが未整備であることにある。がん対策情報センターで公開されている罹患統計は、各年の全国推計値をそのままつなげたものであるため、年によって対象地域や推計手法が異なり解釈が困難である。この問題を解決するために、厚生労働科学研究費補助金第3次対がん総合戦略研究事業の研究班で、地域を限定した年次推移検討手法の開発が進められている。

罹患の年次推移データに関しては、前述の通り、即時性の向上、つまり報告の遅れを解消することがもう一つの課題である。現在約5年遅れになっている罹患統計について、4年程度の予測を行い、死亡統計と同じ年次での最新値の公表を可能にする必要がある。この場合、一定の精度を満たした年次推移データが入手可能な地域が5府県程度に限られているため、時間推計(予測)に加えて、5府県から全国値への空間推計も必要となる。

がんの死亡データであれ罹患データであれ、がんの統計情報は、根拠に基づくがん対策を実施し、評価するための資料である。米国ではACSが、本稿で紹介したような数理的手法を開発し、毎年国および州の統計指標を推計し、対策の評価までを行っている(Jemal et al., 2010)。またACSとは別に、NCIも、がん統計に関する論文を「国への年次報告(Annual Report to the Nation)」として毎年発表している(Edwards et al., 2010)。わが国では、がん対策情報センターががん統計情報を公表しているが、それを科学的に分析、解釈し、価値判断を行う枠組みが未整備である。疫学・数理統計の専門家、政策決定者などが連携して、科学的根拠に基づくがん対策の情報基盤を構築してゆく必要がある。

謝 辞

本稿の作成にあたって貴重な助言をくださった札幌医科大学の加茂憲一先生に心より感謝を申し上げます。

参 考 文 献

- Center for Disease Control and Prevention (2010). Mortality data, <http://www.cdc.gov/nchs/deaths.htm> (Accessed on Dec. 10, 2010).
- Clegg, L. X., Feuer, E. J., Midthune, D. N., Fay, M. P. and Hankey, B. F. (2002). Impact of reporting delay and reporting error on cancer incidence rates and trends, *Journal of National Cancer Institute*, **94**, 1537–1545.
- 独立行政法人国立がん研究センターがん対策情報センター (2010a). 地域がん登録による罹患全国推計の方法, http://ganjoho.jp/professional/statistics/statistics_02.html (Accessed on Dec. 10, 2010).
- 独立行政法人国立がん研究センターがん対策情報センター (2010b). 全国がん罹患モニタリング集計, <http://ganjoho.jp/professional/statistics/monita.html> (Accessed on Dec. 10, 2010).
- 独立行政法人国立がん研究センターがん対策情報センター (2010c). 集計表のダウンロード, <http://ganjoho.jp/professional/statistics/statistics.html> (Accessed on Dec. 10, 2010).
- Edwards, B. K., Ward, E., Kohler, B. A., Ehemann, C., Zaubler, A. G., Anderson, R. N., et al. (2010). Annual report to the nation on the status of cancer, 1975–2006, featuring colorectal cancer trends and impact of interventions (risk factors, screening, and treatment) to reduce future rates, *Cancer*, **116**, 544–573.
- Health Resources and Services Administration, U.S. Department of Health and Human Services (2010). Area resource file, <http://arh.hrsa.gov/> (Accessed on Dec. 10, 2010).
- Jemal, A., Siegel, R., Xu, J. and Ward, E. (2010). Cancer statistics, 2010, *CA Cancer Journal for Clinicians*, **60**, 277–300.
- Kim, H. J., Fay, M. P., Feuer, E. J. and Midthune, D. N. (2000). Permutation tests for joinpoint regression with applications to cancer rates, *Statistics in Medicine*, **19**, 335–351.
- Matsuda, T., Marugame, T., Kamo, K. I., Katanoda, K., Ajiki, W. and Sobue, T. (2011). Cancer incidence and incidence rates in Japan in 2005: Based on data from 12 population-based cancer registries in the monitoring of cancer incidence in Japan (MCIJ) project, *Japanese Journal of Clinical Oncology*, **41**, 139–147.
- National Center for Chronic Disease Prevention and Health Promotion (2010). Behavioral risk factor surveillance system, http://www.cdc.gov/brfss/technical_infodata/index.htm (Accessed on Dec. 10, 2010).
- Pickle, L. W., Feuer, E. J. and Edwards, B. K. (2003). U.S. predicted cancer incidence, 1999: Complete maps by county and state from spatial projection models, NIH Publication No. 03–5435, National Cancer Institute, Bethesda, Maryland.
- Pickle, L. W., Hao, Y., Jemal, A., Zou, Z., Tiwari, R. C., Ward, E., et al. (2007). A new method of estimating United States and state-level cancer incidence counts for the current calendar year, *CA Cancer Journal for Clinicians*, **57**, 30–42.
- 総務省統計局 (2010). 政府統計の総合窓口, <http://www.e-stat.go.jp/SG1/estat/eStatTopPortal.do> (Accessed on Dec. 10, 2010).
- Tiwari, R. C., Ghosh, K., Jemal, A., Hachey, M., Ward, E., Thun, M. J., et al. (2004). A new method of predicting US and state-level cancer mortality counts for the current calendar year, *CA Cancer Journal for Clinicians*, **54**, 30–40.
- Wingo, P. A., Landis, S., Parker, S., Bolden, S. and Heath, B. W. J. (1998). Using cancer registry and vital statistics data to estimate the number of new cancer cases and deaths in the US for the upcoming year, *Journal of Registry Management*, **25**, 43–51.
- World Health Organization (2002). National Cancer Control Programmes: Policies and Managerial Guidelines, 2nd ed., Lyon, France.

Overview of the Mathematical Models Used for Cancer Statistics in the United States of America

Kota Katanoda

Surveillance Division, Center for Cancer Control and Information Services,
National Cancer Center

Cancer statistics are an essential component of the development and evaluation of a national cancer control plan. In Japan, latest cancer mortality data are released after only an approximately one-year delay, while latest cancer incidence data are available after an approximately five-year delay. Thus, updating cancer incidence statistic has been a significant challenge. As of 2010, in the U.S., a state-space model has been used for a three-year prediction of cancer mortality, and a combination of a spatial model, a spatio-temporal model, and a Joinpoint regression model has been used for a four-year prediction of cancer incidence, making both real-time data available. We in Japan also need to prepare timely cancer incidence by developing a short-term prediction method, and to build an information infrastructure for evidence-based cancer control.